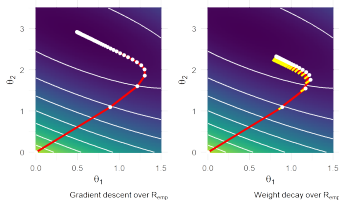


Introduction to Machine Learning

Geometric Analysis of L2 Regularization and Weight Decay



Learning goals

- Have a geometric understanding of $L2$ regularization
- Understand why $L2$ regularization in combination with gradient descent is called weight decay

WEIGHT DECAY VS. L2 REGULARIZATION

Let us optimize the L_2 -regularized risk of a model $f(\mathbf{x} \mid \theta)$

$$\min_{\theta} \mathcal{R}_{\text{reg}}(\theta) = \min_{\theta} \mathcal{R}_{\text{emp}}(\theta) + \frac{\lambda}{2} \|\theta\|_2^2$$

by gradient descent. The gradient is

$$\nabla_{\theta} \mathcal{R}_{\text{reg}}(\theta) = \nabla_{\theta} \mathcal{R}_{\text{emp}}(\theta) + \lambda \theta.$$

We iteratively update θ by step size α times the negative gradient

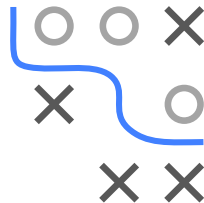
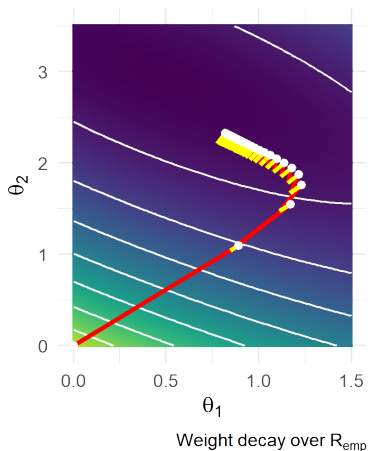
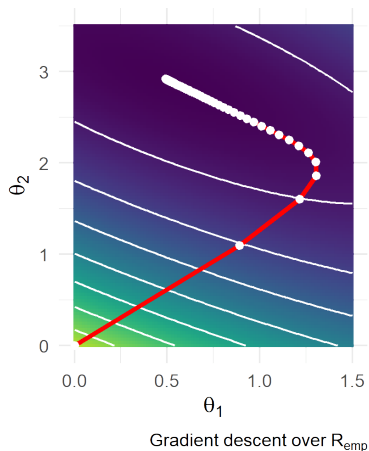
$$\begin{aligned} \theta^{[\text{new}]} &= \theta^{[\text{old}]} - \alpha \left(\nabla_{\theta} \mathcal{R}_{\text{emp}}(\theta^{[\text{old}]}) + \lambda \theta^{[\text{old}]} \right) \\ &= \theta^{[\text{old}]} (1 - \alpha \lambda) - \alpha \nabla_{\theta} \mathcal{R}_{\text{emp}}(\theta^{[\text{old}]}). \end{aligned}$$

The term $\lambda \theta^{[\text{old}]}$ causes the parameter (**weight**) to **decay** in proportion to its size. This is a very well-known technique in deep learning - and simply L_2 regularization in disguise.



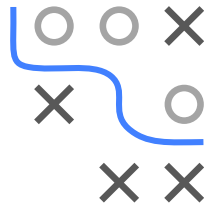
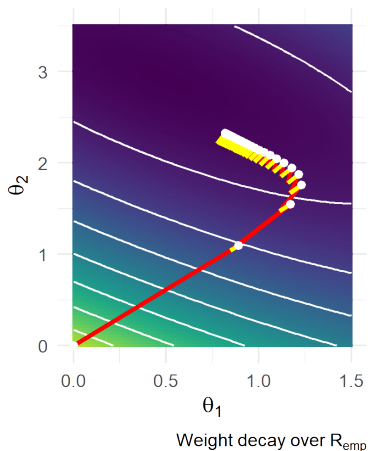
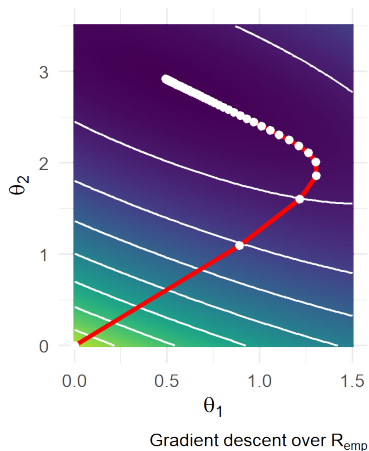
WEIGHT DECAY VS. L2 REGULARIZATION

When we use weight decay, we follow the steepest slope of \mathcal{R}_{emp} as for gradient descent, but in every step, we are pulled back to the origin.



WEIGHT DECAY VS. L2 REGULARIZATION

How strongly we are pulled back to the origin for a fixed stepsize α depends only on λ (as long as the procedure converges):



WEIGHT DECAY VS. L2 REGULARIZATION

Weight decay can be interpreted **geometrically**.

Let's use a quadratic Taylor approximation of the unregularized objective $\mathcal{R}_{\text{emp}}(\boldsymbol{\theta})$ in the neighborhood of its minimizer $\hat{\boldsymbol{\theta}}$,

$$\tilde{\mathcal{R}}_{\text{emp}}(\boldsymbol{\theta}) = \mathcal{R}_{\text{emp}}(\hat{\boldsymbol{\theta}}) + \nabla_{\boldsymbol{\theta}} \mathcal{R}_{\text{emp}}(\hat{\boldsymbol{\theta}}) \cdot (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) + \frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^T \mathbf{H}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}),$$

where \mathbf{H} is the Hessian matrix of $\mathcal{R}_{\text{emp}}(\boldsymbol{\theta})$ evaluated at $\hat{\boldsymbol{\theta}}$.

- The first-order term is 0 in the expression above because the gradient is 0 at the minimizer.
- \mathbf{H} is positive semidefinite.



WEIGHT DECAY VS. L2 REGULARIZATION

The minimum of $\tilde{\mathcal{R}}_{\text{emp}}(\boldsymbol{\theta})$ occurs where $\nabla_{\boldsymbol{\theta}} \tilde{\mathcal{R}}_{\text{emp}}(\boldsymbol{\theta}) = \mathbf{H}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})$ is 0.
Now we $L2$ -regularize $\tilde{\mathcal{R}}_{\text{emp}}(\boldsymbol{\theta})$, such that

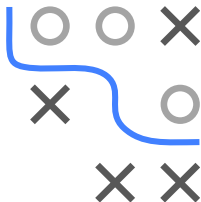
$$\tilde{\mathcal{R}}_{\text{reg}}(\boldsymbol{\theta}) = \tilde{\mathcal{R}}_{\text{emp}}(\boldsymbol{\theta}) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

and solve this approximation of \mathcal{R}_{reg} for the minimizer $\hat{\boldsymbol{\theta}}_{\text{ridge}}$:

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} \tilde{\mathcal{R}}_{\text{reg}}(\boldsymbol{\theta}) &= 0, \\ \lambda \boldsymbol{\theta} + \mathbf{H}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) &= 0, \\ (\mathbf{H} + \lambda \mathbf{I})\boldsymbol{\theta} &= \mathbf{H}\hat{\boldsymbol{\theta}}, \\ \hat{\boldsymbol{\theta}}_{\text{ridge}} &= (\mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}\hat{\boldsymbol{\theta}},\end{aligned}$$

This gives us a formula to see how the minimizer of the $L2$ -regularized version is a transformation of the minimizer of the unpenalized version.

Caveat: Equivalence of weight decay and $L2$ regularization only holds for vanilla SGD (not e.g. Adam)

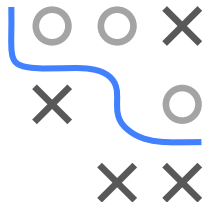


WEIGHT DECAY VS. L2 REGULARIZATION

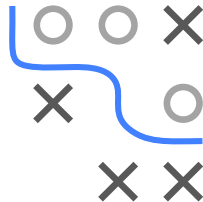
- As λ approaches 0, the regularized solution $\hat{\theta}_{\text{ridge}}$ approaches $\hat{\theta}$. What happens as λ grows?
- Because \mathbf{H} is a real symmetric matrix, it can be decomposed as $\mathbf{H} = \mathbf{Q}\mathbf{\Sigma}\mathbf{Q}^\top$, where $\mathbf{\Sigma}$ is a diagonal matrix of eigenvalues and \mathbf{Q} is an orthonormal basis of eigenvectors.
- Rewriting the transformation formula with this:

$$\begin{aligned}\hat{\theta}_{\text{ridge}} &= \left(\mathbf{Q}\mathbf{\Sigma}\mathbf{Q}^\top + \lambda \mathbf{I} \right)^{-1} \mathbf{Q}\mathbf{\Sigma}\mathbf{Q}^\top \hat{\theta} \\ &= \left[\mathbf{Q}(\mathbf{\Sigma} + \lambda \mathbf{I})\mathbf{Q}^\top \right]^{-1} \mathbf{Q}\mathbf{\Sigma}\mathbf{Q}^\top \hat{\theta} \\ &= \mathbf{Q}(\mathbf{\Sigma} + \lambda \mathbf{I})^{-1} \mathbf{\Sigma}\mathbf{Q}^\top \hat{\theta}\end{aligned}$$

- Therefore, weight decay rescales $\hat{\theta}$ along the axes defined by the eigenvectors of \mathbf{H} . The component of $\hat{\theta}$ that is aligned with the j -th eigenvector of \mathbf{H} is rescaled by a factor of $\frac{\sigma_j}{\sigma_j + \lambda}$, where σ_j is the corresponding eigenvalue.

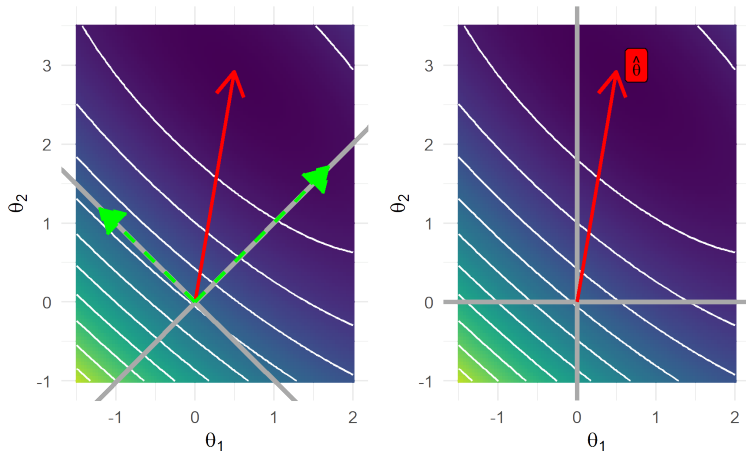


Firstly, $\hat{\theta}$ is rotated by \mathbf{Q}^\top , which we can interpret as a projection of $\hat{\theta}$ on the rotated coordinate system defined by the principal directions of \mathbf{H} :



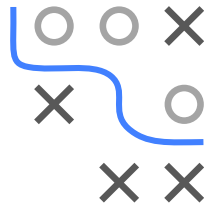
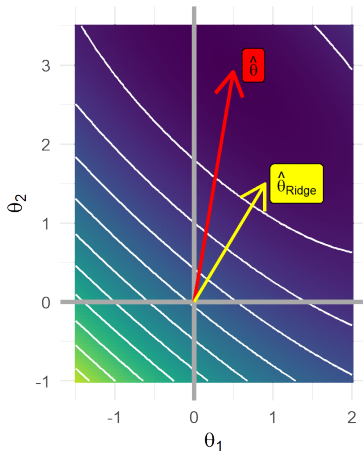
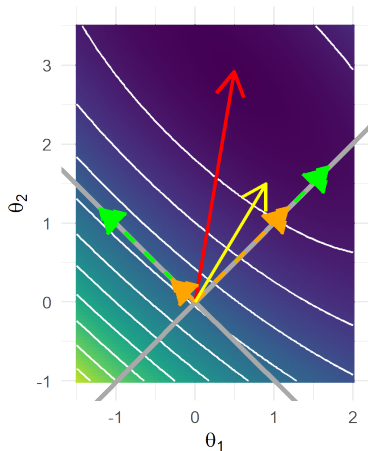
WEIGHT DECAY VS. L2 REGULARIZATION

Since, for $\lambda = 0$, the transformation matrix $(\Sigma + \lambda I)^{-1} \Sigma = \Sigma^{-1} \Sigma = I$, we simply arrive at $\hat{\theta}$ again after projecting back.



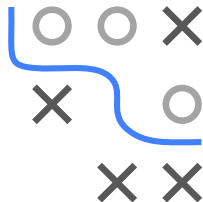
WEIGHT DECAY VS. L2 REGULARIZATION

If $\lambda > 0$, the component projected on the j -th axis gets rescaled by $\frac{\sigma_j}{\sigma_j + \lambda}$ before $\hat{\theta}_{\text{ridge}}$ is rotated back.

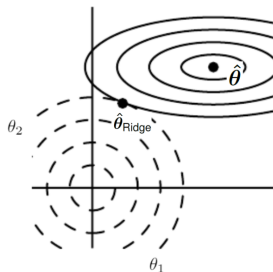


WEIGHT DECAY VS. L2 REGULARIZATION

- Along directions where the eigenvalues of \mathbf{H} are relatively large, for example, where $\sigma_j \gg \lambda$, the effect of regularization is quite small.
- On the other hand, components with $\sigma_j \ll \lambda$ will be shrunk to have nearly zero magnitude.
- In other words, only directions along which the parameters contribute significantly to reducing the objective function are preserved relatively intact.
- In the other directions, a small eigenvalue of the Hessian means that moving in this direction will not significantly increase the gradient. For such unimportant directions, the corresponding components of θ are decayed away.



WEIGHT DECAY VS. L2 REGULARIZATION



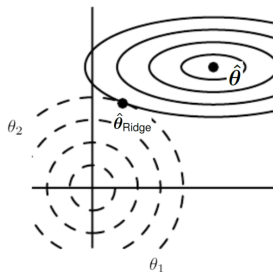
Credit: Goodfellow et al. (2016), ch. 7

Figure: The solid ellipses represent the contours of the unregularized objective and the dashed circles represent the contours of the L_2 penalty. At $\hat{\theta}_{\text{ridge}}$, the competing objectives reach an equilibrium.

In the first dimension, the eigenvalue of the Hessian of $\mathcal{R}_{\text{emp}}(\theta)$ is small. The objective function does not increase much when moving horizontally away from $\hat{\theta}$. Therefore, the regularizer has a strong effect on this axis and θ_1 is pulled close to zero.



WEIGHT DECAY VS. L2 REGULARIZATION



Credit: Goodfellow et al. (2016), ch. 7

Figure: The solid ellipses represent the contours of the unregularized objective and the dashed circles represent the contours of the L_2 penalty. At $\hat{\theta}_{\text{ridge}}$, the competing objectives reach an equilibrium.

In the second dimension, the corresponding eigenvalue is large indicating high curvature. The objective function is very sensitive to movement along this axis and, as a result, the position of θ_2 is less affected by the regularization.

